

# Extending ChatGPT with a Browserless System for Web Product Price Extraction

Jorge Lloret-Gazo

Dpto. de Informática e Ingeniería de Sistemas.  
Facultad de Ciencias. Edificio de Matemáticas.  
Universidad de Zaragoza. 50009 Zaragoza. Spain.  
jlloret@unizar.es

**Abstract.** With the advenement of ChatGPT, we can find very clean, precise answers to a varied amount of questions. However, for questions such as 'find the price of the lemon cake at zingerman's', the answer looks like 'I can't browse the web right now'. In this paper, we propose a system, called Wextractor, which extends ChatGPT to answer questions as the one mentioned before. Obviously, our system cannot be labeled as 'artificial intelligence'. Simply, it offers to cover a kind of transactional search that is not included in the current version of ChatGPT. Moreover, Wextractor includes two improvements with respect to the initial version: social extraction and pointing pattern extraction to improve the answer speed.

## 1 Introduction

Since its launch, ChatGPT has grown quickly to become a well-known question-and-answer system. Numerous exam-related use cases involving ChatGPT have had their performance examined, with different degrees of scientific rigor ranging from in-depth research to anecdotal evidence. Use cases include excellent success on the United States Medical Licensing Examination [8]. For these and other reasons, large language models (LLMs) are expected to have an effect on a variety of fields and be used as assistants by a variety of professionals .

However, there are many easier questions that ChatGPT is not able to answer, for example, those related to queries that requires browser capabilities and, in particular, with product prices on the web.

On the other hand, businesses are interested in extracting product prices from different web sites and in following them up during a certain period of time for different purposes. For example, an e-commerce enterprise is interested in the price of pairs of shoes listed on different competing websites with the purpose of adapting their prices when competence prices fall below a certain threshold.

In this paper, we propose the combination of ChatGPT with an enhanced version of Wextractor [12] to answer questions about product prices.

To be specific, in paper [12], we proposed an architecture for extracting prices that applies the following four steps: (0) browserless extraction of raw HTML, (1) fragmentation, (2) rule application and (3) price extraction. In step (0), raw

HTML is obtained from the url of the page. In step (1), fragments that contain clues of the prices are found. In the next step, rules previously designed by a rule designer are applied to the fragments. As a result, several fragments, that contain several candidate prices, remain and one of them is chosen as the target value. Step 0 is based on the browserless concept of [4]. The first step is inspired by the segmentation of [19] while, to the best of our knowledge, there is no proposal of easy-to-specify rules in the literature for the purpose of extracting prices.

In this paper, we extend the original architecture with two ideas: (a) social extraction and (b) pointing pattern extraction. The first idea consists of storing prices of objects coming from searches of other users and using these prices to answer subsequent queries that occur near the initial search. The second idea is to use automatic pointing pattern extraction. A pointing pattern is a regular expression that points to where the entity price is.

Then, the second and subsequent times an extraction is performed on a page, first, a social search will be performed. If the price is available in the social storage, it will be returned as a result. Otherwise, the pointing pattern of the page is used. If this extraction does not find any price, the process begins and steps (1)-(3) are applied again (step (0) was previously applied, so it is not applied again).

Taking into account our enhanced architecture, questions on product prices can be dispatched from ChatGPT to Wextractor, extending in this way the capabilities of ChatGPT.

The contribution of this paper is threefold:

1. We offer an answer to ChatGPT-uncovered questions
2. We incorporate a social module into Wextractor in which people act as sensors that transmit information from the market to our architecture. Therefore, some searches can be answered with the information previously gathered from users. This avoids computationally expensive queries.
3. Finally, we incorporate a pointing pattern module so that a faster search can be performed using a pointing pattern.

The rest of the paper is organised as follows. Section 2 is devoted to defining the problem and the challenges. In Section 3, we explain the architecture and in Section 5, an implementation of it. Section 4 is devoted to the connexion between ChatGPT and Wextractor. In the last sections, we deal with the evaluation, related work, conclusions and future work.

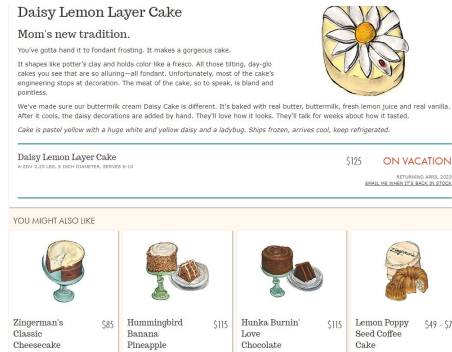
## 2 Problem definition and description of the architecture

Although ChatGPT is a powerful tool, there are questions which remain not satisfactorily answered. For example, let us consider a query such as ‘price of lemon cake at Zingerman’s’. The answer of ChatGPT is shown in Figure 1. That is, ChatGPT cannot deal with real-time queries or with queries that need browser capabilities.

find the price of the lemon cake at zingerman's

I can't browse the web right now. Can I help you with something else? You might be able to find the price of the lemon cake by checking Zingerman's official website or contacting their customer service.

**Fig. 1.** Answer of ChatGPT to a query about Zingerman's cake



**Fig. 2.** Fragment of the Zingerman web page as of March 2024

On the other hand, our Wextractor architecture responds to questions as follows: Given a detail web page  $pg$  of an entity  $e$  with url  $u$ , extract the price of  $e$  from page  $pg$ .

An example of this problem is query ( $q_1$ ): Given the url

<https://www.zingermans.com/Product.aspx?ProductID=A-ZDV>

find the price of the lemon cake.

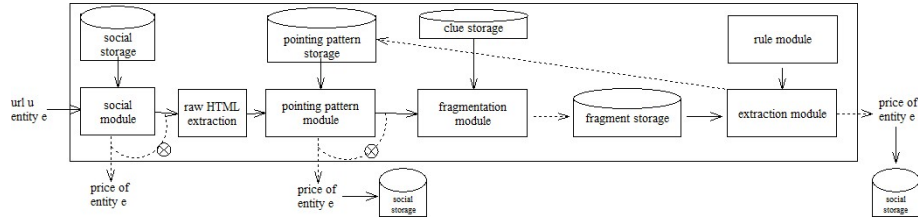
The problems we solve in this paper are:

1. How to improve the answer speed of Wextractor (see Section 6).
2. How to connect ChatGPT with Wextractor (see Section 4) to answer real-time queries on prices with ChatGPT.

The solution of the first problem is expressed in the following four steps. As an example, we will use query ( $q_1$ ). Together with the steps, we summarize our architecture extended with the social and the pointing pattern modules (see 3). Additional details of the original architecture can be found on [12].

### Step 0. Social extraction

This extraction is done by the *social module*. It receives the url  $u$  of a detail page for an entity  $e$  and searches the price associated with the url  $u$  in the *social storage*. If the price is found, it is returned provided that the difference of the actual time with the time of the price in the social storage falls within a limit. If the URL is not available or the limit is overcome, then jump to the step 1.



**Fig. 3.** Our architecture

This module is inspired by Waze [18]. In its page, we can read “Nothing can beat people working together”. The idea behind this module is that people act as sensors that transmit the information from the market to the Wextractor architecture. This information gathered by the sensors is elaborated inside our architecture and offered curated to other people doing similar searches through ChatGPT.

In this way, information is bidirectional. As far as I know, the information flows from the buyer to the seller and remains stored there or is stored in Google. According to our idea, other buyers can directly access the information of other buyers.

### Step 1. Browserless extraction of raw HTML

From the url  $u$  of the page, the raw HTML code is extracted without interaction from a browser.

### Step 2. Pointing pattern extraction

This extraction occurs when the social extraction does not give any results. It needs the following two modules:

*Pointing pattern storage* This component stores pairs formed by the url of a page and a pointing pattern, i.e., a regular expression that points to the price of the page.

*Pointing pattern module* This module receives the url of the page of an entity and searches in the *pointing pattern storage* for any pointing pattern associated with the url  $u$ . If a pointing pattern is found, the price is extracted from the raw HTML of the page obtained in Step 1.

As a result, the price of the entity or nothing if the pointing pattern does not work is obtained. In this last case or if the page with url  $u$  was not previously accessed, an extraction from scratch begins in the next step.

**Step 3. From scratch extraction** From this point on, a from scratch extraction is performed. This extraction occurs in one of the following cases:

1. it is the first time a page is accessed
2. the page was previously accessed but neither the social extraction nor the pointing pattern extraction gave any results

It needs the following modules:

*Clue storage* This stores strings representing the currencies of interest. For example, EUR for the currency euro.

*Fragmentation module* This module extracts, from the raw HTML of a page with url  $u$ , the fragments that contain possible prices of entity  $e$ . To do so, it uses the *clue storage* to identify the places where prices are. Examples of clues are  $\text{\&euro;}$ ; or  $\text{\$}$ . For query (q1), 149 fragments were extracted, two of them being:

(f1) `<span id="ct100" class="price">\$ 125</span>`

(f2) `<div class="saving">SAVE10%=\&euro;;12.80</div>`

All these fragments are stored in the *fragment storage*.

*Rule module* The rule module is composed of an ordered set of Condition-Action rules elaborated by the designer. We will use *discarding rules* to discard fragments verifying conditions that indicate the absence of the target price. Then, a weight of one is added to these fragments and zero to the rest of the fragments. In the end, fragments with a total weight equal to zero are selected.

In the example of the cake, the discarding rules search for properties of the fragments that indicate that they do not contain the price of the cake. An example of a discarding rule, called `semr1`, is ‘discard those fragments that contain the word SAVE’. (f2) is an example of a fragment discarded by the `semr1` discarding rule.

*Extraction module* This module applies the rules of the *rule module* to the fragments of the *fragment storage* and determines the right price for entity  $e$ . This module works as follows. First, the fragments of the page are loaded from the *fragment storage* and their weight is initialized to zero. Next, the rules are applied to the fragments updating the weights accordingly. Then, the number of candidate prices is calculated where the candidate price are those whose fragments have weight equal to zero. If only one candidate price is found this is the target price.

In the example of the cake, fragment (f1) is the only non-discarded fragment and from it the value  $\text{\$}125$  is obtained.

However, no price or more than one price could also be retrieved. To remedy this, we have defined recovering rules, i.e., rules that try to rescue the correct price from previously discarded prices. After applying these rules, if only one price is found this is the target price. Otherwise, there is no solution.

If the user wants to check the price of an entity again, we have two options. One of them is to repeat a from scratch extraction. The second one is to try to use the information of the from scratch extraction to speed up the process. Because our method relies on detecting HTML fragments and, in general, the web page structure does not change too much with time, we have chosen the second option.

Following the second option, the non-discarded fragments of the from scratch extraction are used to build a pointing pattern, that is, a regular expression that matches the price. For example, after applying the from scratch extraction in query (q1), only fragment (f1) remains. From this fragment, the pointing pattern `pp1 id="ct100" class="price">\$[0-9]{2,3}\.[0-9]{1,2}` is built for subsequent extractions. It is worth noting that when a new extraction is performed with this pattern, the price of  $\text{\$}125$  could have fallen below  $\text{\$}100$ . For this reason, the

first part of the numeric pattern is  $[0-9]\{2,3\}$ , that is, prices of less than 100\$ are also matched by the pattern.

According to the previous discussion, two additional actions are taken after the from scratch extraction success:

1. A pointing pattern is extracted and stored in the *pointing pattern storage* to be used in future searches
2. Quaterns (e,u, p, t) are stored in the *social storage* where e is an entity, u is the url of a detail page for entity e in page u and t is the timestamp where the price was found

Then, the next time ChatGPT demands the price of an entity e, it will dispatch the query to the enhanced Wextractor architecture which, in turn, will offer an answer, such as the url of the page together with the price of 125 \$. Otherwise, the query would remain without an answer inside ChatGPT.

### 3 Combined architecture with ChatGPT and Wextractor

In this Section, we explain how to combine the potential of ChatGPT and Wextractor to answer price questions.

The way of working is as follows (see Figure 4): The user issues a query against ChatGPT which, extracts the tuple (e0,s0) from the query, where e0 is an entity for which the price must be found in site s0. For example, from the query ‘find the price of the lemon cake at Zingerman’, ChatGPT extracts the pair (‘lemon cake’,’Zingerman’).

On the other hand, the Wextractor architecture is fed with user searches for prices. In these searches, the user provides a url from which a price is obtained. However, this is not the case for ChatGPT where the user provides a textual query, without URLs, such as ‘find the price of of the lemon cake at Zingerman’. So, in order to communicate ChatGPT and Wextractor, we need a new module to find the urls of Wextractor that correspond to the site of the ChatGPT query.

This new module is called `siteTOURLMatcher` and is responsible for finding every url in the social storage that belongs to the site (see Figure 4). As a result, the module obtains a pair (possibly empty) (u,p) where u is an url corresponding to the site s0 and p is the price of e0 in the page with url u.

For example, for the pair (‘lemon cake’,’Zingerman’), the `siteTOURLMatcher` module queries the social storage of Wextractor and finds the pair

(<https://www.zingermans.com/Product.aspx?ProductID=A-ZDV>, \$125)

So, this pair is sent back to ChatGPT as an answer to the query ‘find the price of the lemon cake at Zingerman’.

### 4 An implementation of the Architecture

In this section, we describe a specific implementation of our architecture. As designers, we have chosen to implement our architecture by using relational databases and SQL, all glued with a procedural language for relational databases.

**Table 1.** Pseudocode for finding the price of a detail page

Algorithm Wextractor

---

Input: url  $u$  of the page  $pg$  of an entity  $e$

Output:

success=-1. The page is not available

success=0. No price was found.

success=1. One price was found by from scratch.

success=2. One price was found by pointing pattern extraction.

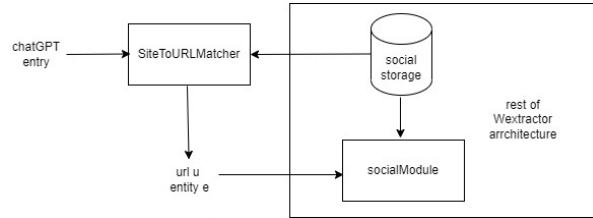
success=3. One price was found by social extraction.

price of the entity when success=1 or 2 or 3

---

Pseudocode

```
1. success<-false
2. if isAvailablePage(u) then
3.   found=getSocialPrice(u, price, timestamp)
4.   if(found) then
5.     if(sysdate-timestamp<socialExtractionValidity) then success=3
6.     else found=false
7.   end if
8.   end if
9.   if (!found)
10.    found=getPointingPattern(pointingPattern, u)
11.    if (found) then
12.      if(sysdate-timestamp<pointingPatternValidity) then
13.        getHTMLCode(HTMLcode, u)
14.        found=doPPEExtraction(HTMLCode, pointingPattern, price)
15.        if (found) then success=2
16.        else found=false
17.      end if
18.    else found=false
19.  end if
20.  end if
21.  end if
22.  if !found
23.    getHTMLCode(HTMLcode, u)
24.    found=doFromScratchExtraction(HTMLCode,s)
25.    if (found) then success=1
26.    findAndStorePointingPattern(HTMLCode, u)
27.    storageInSocialStorage(entity, u,price, timestamp)
28.  else success=0
29.  end if
30.  end if
31. else
32. success=-1
33. end if
```



**Fig. 4.** Connection between ChatGPT and Wextractor

Table 1 shows the main algorithm and combines the three modules of the architecture.

We have set a time of validity of a social extraction (`socialExtractionValidity`) and a time of validity of a pointing pattern extraction (`pointingPatternValidity`) where `socialExtractionValidity < pointingPatternValidity`. The

`socialExtractionValidity` means that a correct price obtained at time `t0` is valid for any other request until time `t0+socialExtractionValidity`. If this limit is overcome, the price is extracted by a pointing pattern search until `t0+pointingPatternValidity`. We consider that social extraction is less durable because it deals with the prices of the object and pointing pattern is more durable because it deals with the structure of the page, which changes less frequently.

In the algorithm, a social extraction is first performed (lines 3 to 8). If this extraction does not give the expected results, a pointing pattern extraction is performed (lines 9 to 21). In particular, the first time an extraction is performed, as there are no previously available patterns, a from scratch extraction is always performed (lines 22 to 30).

In detail, the algorithm is as follows: In line 2, the existence of a page with url `u` is checked. If the page is not available, `success=-1` is returned (line 32).

**Social extraction** If the page is available in the social storage (line 3) and the time of the social storage falls within a limit, then `success=3` (line 5) and the price fetched from the social storage is the result.

**Pointing pattern extraction** Otherwise, in line 10, a pointing pattern is searched in the pointing pattern storage for page `u`. If it is found and the time of the pointing pattern falls within a limit, a pointing pattern extraction is performed (line 14), using the algorithm `doPPEExtraction`. At the end, if a price is found, it is the result and `success` is set to 2 (line 15). Otherwise, the pointing pattern extraction failed because, for example, the structure of the page has changed so that the pointing pattern does not find the correct price.

**From scratch extraction** If the previous extractions did not find any solution, a from scratch extraction is performed (line 24). If a price is found, it is the result and `success` is set to 1. Moreover, a pointing pattern is determined (line 26) and added to the pointing pattern storage for use in subsequent pattern-based extractions. On the other hand, the entity, url, found price and the timestamp in which the price was found are stored in the social storage (line 27). If no price is found, `success` is set to 0 (line 28).



## 5 Experimental validation

Our architecture could be tested in a real environment, but it would take several months to complete. Instead, we prepared a simulation based on the data gathered from our 735-page dataset [12]. To be specific, our previous dataset served to set the number of pages and the probability of success. These initial data have been combined with our hypothesis that the request of pages follows a Zipf distribution. Let us dive into the details.

**Dataset** We searched for websites where we could find very visited shopping websites with the following requirements: (a) the list must be representative of the most popular shopping websites globally, and (b) the list must consist of shopping websites in English so that we would have the means to analyze the data collected from the websites.

Previously, we could retrieve the list from places like Alexa using the Top Sites API [9] or WebShrinker or 5000best.com. However, at the time of writing, Alexa is not available, WebShrinker does not have an open version of the data and 5000best.com does not offer any information about shopping websites. Moreover, a Google search on the topic gives very generic results. For these reasons, we will use our original list of 735 pages published in the paper [12].

**Simulation Features** The general features of our simulation are as follows:

1. The simulation runs on  $N_{\text{pages}}$  web pages from dataset D.
2.  $N \geq 0$  distinct requests are made of dataset D in the integer values of the time interval  $[1, N]$ , that is, a request per unit of time
3. Each request has the form  $(p, t)$ , where  $p$  belongs to dataset D and  $t$  is an integer time of interval  $[1, N]$
4. The pages are demanded according to a Zipf distribution of parameters  $a=b=1$ . Using the algorithm [3], we generated  $N$  numbers between 1 and  $N_{\text{pages}}$  each simulating a requested page. Without loss of generality, we supposed that the pages are ordered according to their demand. That is, page 1 is the most demanded and page  $N_{\text{pages}}$  is the least demanded. Therefore, page 2 is half demanded than page 1.
5. The probability of success (that is, the price was found) of the from scratch search on each page is  $p$  and the probability of failure is  $1-p$ .

### Simulation at Work

**Exact numbers** For simulating our architecture at work, we used the following figures:  $N_{\text{pages}}=735$ ,  $N=50.000$ ,  $p=579/735$ ,  $\text{socialExtractionValidity}=10$ ,  $\text{pointingPatternValidity}=20$ .

For determining  $p$ , we used the fact that our Wextractor architecture produces exact results for 579 out of 735 pages [12]. Therefore, we calculated the probability of success for each page in a from scratch extraction as  $p=579/735$ .

**Simulation database** We implemented the simulation using the Oracle database 21g and the PL/SQL programming language. The database upon which the simulation is built consists of the following tables.

```
fromScratchResultOfExtraction(idPage, success)
```

**Table 2.** Frequency of success

Page	Success=1	Success=2	Success=3
1	412	1382	5073

```
zipfNumber(id, numero)
resultOfExperiment(idPage, success, time)
```

The data input are in tables `fromScratchResultOfExtraction` and `zipfNumber`. The table `fromScratchResultOfExtraction` has been filled with 735 rows, one per page, and the success column is calculated according to the probability  $p=579/735$ . The table `zipfNumber` is filled with  $N$  values of `idPage` where each value ranges between 1 and `Npages`. The values follow a Zipf distribution, as stated in point 4.

The output data are in table `resultOfExperiment`. This table is completed after executing the simulation algorithm below under the assumption that the results of the from scratch extraction are in table `fromScratchResultOfExtraction`. The table `resultOfExperiment` will contain  $N$  rows and each row is composed of the id of the queried page, the success of the search and the time of the search. The possible values for success are: 0 if no price is found, 1 if the price is found from scratch, 2 if the price is found by pointing pattern extraction and 3 if the price is found by social extraction.

**Simulation algorithm** With this database, the simulation algorithm is as follows

```
for each i in 1..N
  page<-getZipfNumber(i)
  success<- getSuccess(page, i)
  storeSuccess(page, success)
end for
```

**Results and conclusions** The number of successful experiments were 43130 out of 50.000. This means 86.26% success.

With respect to individual pages, in Table 2, we gather the frequency of success for the most frequent page, number 1. Page number 1 was requested 6867 of 50.000 times. As can be seen in Table 2, in most cases it was not needed a from scratch extraction, leading to the quickest answer by the social extraction or pointing pattern extraction.

Finally, for those frequent pages for which our Wextractor does not find any solution, it would be better to manually find a pointing pattern. If we found manually a pointing pattern for the 10% most frequent pages, the success rate would increase from 86.26% to 90.67%.

**Reproducibility** Following the reproducibility recommendation of the ACM SIGMOD and of VLDB, we have made available at [1] the script for an implementation of the architecture, the text files with the HTML code of the web pages of our dataset and the code of the simulation of how the architecture works.

## 6 Related Work

Recently, ChatGPT has emerged as a question-and-answer system. However, as evidenced in this paper, certain inquiries, like web prices, remain challenging to address effectively. To our awareness, no existing research has extended ChatGPT to encompass the capability of extracting prices from the web, unlike our present work. This paper aims to bridge this gap by introducing an innovative approach to augmenting ChatGPT with price extraction capabilities. By doing so, we contribute to enhancing ChatGPT’s utility in addressing a broader range of inquiries, particularly those pertaining to e-commerce and data mining. Our endeavor underscores the importance of continually advancing AI systems to tackle increasingly complex tasks and highlights the potential for extending ChatGPT’s functionality beyond its current scope. Through this research, we aim to facilitate more comprehensive and insightful interactions with ChatGPT, thereby enriching its practical applications in various domains.

With respect to the problem of extracting data from web pages, it appeared early in the year 2000 and it is an important task in e-commerce and data mining. The proposed solution has essentially been the use of wrappers, that is, are programs or scripts that help extract structured data from web pages by identifying specific HTML elements or patterns. Reviewing the literature, we have three main ways of generating wrappers [2, 5, 15]: manual, wrapper induction, visual based.

**Manual creation.** The manual creation of wrappers [17] is a traditional approach used for extracting data from web pages. It involves manually designing and coding scripts or programs to extract specific information from HTML documents. This method requires a deep understanding of web page structures and the ability to write custom code to navigate and extract data accurately. In the manual creation process, developers typically analyze the target web pages, identify the relevant HTML elements, and define the extraction rules. These rules can include techniques such as pattern matching, XPath queries, or CSS selectors to locate and extract the desired data, such as prices. The extracted data is then processed and transformed for further analysis or integration into other systems.

One advantage of manual wrapper creation is its flexibility and control over the extraction process. Developers can customize the extraction logic to handle various website layouts and adapt to changes over time. It also allows for fine-grained control over data quality and validation. However, manual wrapper creation can be time-consuming and labor-intensive, especially when dealing with a large number of websites or complex web page structures. It requires technical expertise and ongoing maintenance to keep the wrappers up-to-date as websites evolve.

To streamline the process, researchers have explored semi-automatic or automatic methods for wrapper creation, such as machine learning algorithms or rule induction techniques. These approaches aim to reduce the manual effort required by automatically learning extraction rules from training data or by inferring patterns from example pages.

**Wrapper induction creation** Wrapper induction [9,10,16], also known as automatic wrapper generation, is an approach that aims to automate the process of creating wrappers for data extraction from web pages. Instead of manually designing and coding extraction rules, wrapper induction techniques use machine learning algorithms to automatically learn the extraction patterns from example web pages.

In the wrapper induction process, a set of example web pages with the desired data is provided as training data. The machine learning algorithm analyzes the HTML structures, tags, attributes, and textual content of these pages to identify patterns and create extraction rules. These rules are then used to extract the desired data from unseen web pages that have a similar structure.

One of the key advantages of wrapper induction is its ability to handle a large number of web pages and adapt to changes in website layouts. As the algorithm learns from the training examples, it generalizes the extraction rules to handle variations in web page structures and accurately extract the desired data. This makes it a scalable and efficient approach for data extraction from diverse sources.

Wrapper induction techniques can employ various machine learning algorithms, such as decision trees, rule-based systems, or sequence models, depending on the characteristics of the data and the extraction task. Some approaches also incorporate active learning, where the algorithm iteratively selects informative examples for labeling by a human expert, further improving the accuracy of the wrapper generation process.

However, wrapper induction may face challenges when dealing with noisy or complex web pages that deviate from the patterns observed in the training data. Handling dynamic content, JavaScript-driven interactions, or dynamically generated elements can also be challenging for wrapper induction algorithms.

**Visual based wrapper creation** Visual-based wrapper creation (Mozenda [13], iMacros [7], Visual Web Ripper [14], Lixto [11]), also known as visual web scraping or visual data extraction, is an approach that leverages visual information to extract data from web pages. Unlike manual or programmatic methods that rely on HTML structures and coding, visual-based approaches utilize computer vision and machine learning techniques to recognize and extract data from the visual elements of web pages.

In visual-based wrapper creation, a user interacts with a tool or framework that provides a visual interface to define the extraction process. The user typically highlights the desired data elements on a web page, and the system automatically generates the extraction rules based on the visual cues. These extraction rules can include spatial relationships, visual patterns, or templates to identify and extract the relevant data.

One advantage of visual-based wrapper creation is its accessibility to users with limited programming knowledge. It allows non-technical users to easily define extraction rules by visually selecting the data elements of interest. This makes the process more intuitive and user-friendly, enabling a wider range of individuals to extract data from web pages without extensive coding expertise.

Visual-based wrapper creation also facilitates adaptability to changes in web page layouts. Since it relies on visual patterns rather than HTML structures, it can handle variations in web page designs more effectively. When a web page undergoes updates or modifications, the visual-based approach can adapt by identifying the visually similar elements for data extraction.

While visual-based wrapper creation offers simplicity and flexibility, it may face challenges with complex or dynamic web pages that contain intricate visual elements or require interaction to reveal the desired data. Additionally, the accuracy of the extraction heavily relies on the quality of the visual recognition algorithms and the ability to handle variations in visual presentation.

**Comparison with our work** We have found no published studies papers proposing a method similar to ours. However, browserless extraction is based on a previous paper [4]. Our segmentation phase is inspired by the segmentation described in [19]; however, to the best of our knowledge, there is no proposal reported in the literature of discarding rules combined with social extraction and pointing pattern extraction. Finally, most of studies attempt to extract all the information from a page so that they present the price of the main entity of interest as well as the price of every secondary entity that appears on the page. They do not include procedures to discriminate between the main price and the secondary prices as we do with our approach. See, for example, [6].

## 7 Conclusions and Future work

In this study, we successfully augmented ChatGPT with price extraction functionality by implementing our Wextractor architecture. This integration was necessitated by ChatGPT's inherent limitation in addressing queries requiring browsing capabilities. Our approach represents a loosely coupled integration, as Wextractor operates independently of ChatGPT's artificial intelligence capabilities.

In addition, enhancements were made to our architecture to optimize response time, particularly through improvements in social extraction and pointing pattern extraction techniques. These optimizations contribute to a more efficient and timely retrieval of information.

Looking ahead, a critical focus of future research involves extending ChatGPT's capabilities to handle real-time queries requiring browser functionalities. This presents a significant challenge that requires careful consideration. One avenue for exploration is the potential reconstruction of Wextractor from the ground up, aligning its functionalities more closely with the requirements of a large language model such as ChatGPT. Alternatively, we must explore methods to integrate the innovative ideas and functionalities of Wextractor into the existing ChatGPT framework.

In conclusion, our study represents a significant step forward in augmenting ChatGPT's capabilities for practical applications. However, the journey does not end here; rather, it opens up avenues for further research and development. By addressing the challenges of real-time browsing capabilities and exploring inte-

gration strategies, we can continue to enhance the effectiveness and versatility of ChatGPT in addressing several user inquiries and tasks.

## References

1. Reproducibility of our price extraction architecture. <https://acortar.link/96klmh>. Accessed: 2024-02-26.
2. Chia-Hui Chang, Mohammed Kayed, Moheb R Girgis, and Khaled F Shaalan. A survey of web information extraction systems. *IEEE transactions on knowledge and data engineering*, 18(10):1411–1428, 2006.
3. Kenneth J. Christensen. <https://cse.usf.edu/~kchrise/tools/genzipf.c>.
4. Ruslan R. Fayzrakhmanov, Emanuel Sallinger, Ben Spencer, Tim Furche, and Georg Gottlob. Browserless web data extraction: Challenges and opportunities. In *Proceedings of the 2018 World Wide Web Conference, WWW '18*, pages 1095–1104, Republic and Canton of Geneva, Switzerland, 2018. International World Wide Web Conferences Steering Committee.
5. Emilio Ferrara, Pasquale De Meo, Giacomo Fiumara, and Robert Baumgartner. Web data extraction, applications and techniques: A survey. *Knowledge-based systems*, 70:301–323, 2014.
6. Tim Furche, Georg Gottlob, Giovanni Grasso, Christian Schallhart, and Andrew Jon Sellers. Oxpath: A language for scalable data extraction, automation, and crawling on the deep web. *VLDB J.*, 22(1):47–72, 2013.
7. iMacros. <http://www.iopus.com/imacros/>.
8. Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198, 2023.
9. Nicholas Kushmerick, Daniel S. Weld, and Robert B. Doorenbos. Wrapper induction for information extraction. In *IJCAI (1)*, pages 729–737, 1997.
10. Alberto H. F. Laender, Berthier A. Ribeiro-Neto, and Altigran Soares da Silva. Debye - data extraction by example. *Data Knowl. Eng.*, 40(2):121–154, 2002.
11. Lixto. <http://www.lixtto.com>.
12. Jorge Lloret-Gazo. A browserless architecture for extracting web prices. In *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, pages 2193–2200, 2020.
13. Mozenda. <http://www.mozenda.com>.
14. Visual Web Ripper. <http://www.visualwebripper.com/>.
15. Sunita Sarawagi. Information extraction. *Foundations and Trends in Databases*, 1(3):261–377, 2008.
16. Stephen Soderland. Learning information extraction rules for semi-structured and free text. *Machine Learning*, 34(1-3):233–272, 1999.
17. Jakub Starka, Irena Holubova, and Martin Necasky. Strigil: A framework for data extraction in semi-structured web documents. In *iiWAS*, page 453, 2013.
18. Shoshana Vasserman, Michal Feldman, and Avinatan Hassidim. Implementing the wisdom of waze. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
19. Yanhong Zhai and Bing Liu. Structured data extraction from the web based on partial tree alignment. *IEEE transactions on knowledge and data engineering*, 18(10):1614–1628, 2006.